*[Continued on next page]*

(54) Title: CLUSTERING OF TEXT FOR STRUCTURING OF TEXT DOCUMENTS AND TRAINING OF LANGUAGE MODELS

(57) Abstract: The present invention relates to a method, a text segmentation system and a computer program product for clustering of text into text clusters representing a distinct semantic meaning. The text clustering method identifies text portions and assigns text portions to different clusters in such a way that each text cluster refers to one or several semantic topics. The clustering method incorporates an optimization procedure based on a re-clustering procedure evaluating a target function being indicative of the correlation between a text unit and a cluster. The text clustering method makes use of a text emission model and a cluster transition model and makes further use of various smoothing techniques.

# WO 2005/050473 A2

For two-letter codes and other abbreviations, refer to the "Guid-ance Notes on Codes and Abbreviations" appearing at the begin-ning of each regular issue of the PCT Gazette.